

# Codes with Locality for Two Erasures

N. Prakash, V. Lalitha and P. Vijay Kumar

## Abstract

In this paper, we study codes with locality that can recover from two erasures via a sequence of two local, parity-check computations. By a local parity-check computation, we mean recovery via a single parity-check equation associated to small Hamming weight. Earlier approaches considered recovery in parallel; the sequential approach allows us to potentially construct codes with improved minimum distance. These codes, which we refer to as locally 2-reconstructible codes, are a natural generalization along one direction, of codes with all-symbol locality introduced by Gopalan *et al*, in which recovery from a single erasure is considered. By studying the Generalized Hamming Weights of the dual code, we derive upper bounds on the minimum distance of locally 2-reconstructible codes and provide constructions for a family of codes based on Turán graphs, that are optimal with respect to this bound. The minimum distance bound derived here is universal in the sense that no code which permits all-symbol local recovery from 2 erasures can have larger minimum distance regardless of approach adopted. Our approach also leads to a new bound on the minimum distance of codes with all-symbol locality for the single-erasure case.

## I. INTRODUCTION

A primary goal in distributed data storage is the efficient repair of a failed node. While regenerating codes [1] aim to minimize the amount of data download needed to carry out node repair, codes with locality [2] seek to minimize the number of nodes accessed during node repair. The focus of the present paper is on codes with locality.

Let  $\mathcal{C}$  denote an  $[n, k, d_{\min}]$  linear code having block length  $n$ , dimension  $k$  and minimum distance  $d_{\min}$ . Where the minimum distance is not relevant, we will simply refer to  $\mathcal{C}$  as an  $[n, k]$  code. The  $i^{\text{th}}$  code-symbol  $c_i$ ,  $1 \leq i \leq n$ , of the code  $\mathcal{C}$  is said to have locality  $r$  if this symbol can be recovered by accessing at most  $r$  other code symbols and performing a linear computation. Equivalently, there exists a row in the parity-check matrix  $H$  of the code of Hamming weight  $\leq (r+1)$ , whose support includes  $i$ . A systematic code in which all the  $k$  message symbols have locality  $r$  is said to have information locality  $r$ . The minimum distance  $d_{\min}$  of a code with information locality  $r$  is upper bounded [2] by

$$d_{\min} \leq n - k - \left\lceil \frac{k}{r} \right\rceil + 2. \quad (1)$$

The pyramid-code construction in [3] yields optimal codes for all  $\{n, k, r\}$  with field size  $O(n)$ . The authors of [2] also introduce the notion of all-symbol locality in which all code symbols, not just the message symbols, have locality  $r$ . They show the existence of codes with all-symbol locality that achieve the bound in (1) when  $(r+1) \mid n$ , but leave open the question as to whether it is possible to derive a tighter bound in the all-symbol-locality case, for general  $\{n, k, r\}$ . The all-symbol-locality property is preferable in applications as it permits a uniform approach to storage-system design. Codes with locality also go by the name locally-repairable [4], locally-reconstructible [5] and locally-recoverable codes [6].

### A. Handling Multiple Erasures

There is current practical interest in the handling of multiple erasures as simultaneous node failures are not uncommon, given the increasing trend towards replacing expensive servers with low-cost commodity servers, the presence of “hot” nodes etc. Several approaches to the multiple-erasure case in the context of codes with locality can be found in the literature.

N. Prakash, V. Lalitha and P. Vijay Kumar are with the Department of ECE, Indian Institute of Science, Bangalore, 560 012 India (email: {prakashn, lalitha, vijay}@ece.iisc.ernet.in).

This research is supported in part by the National Science Foundation under Grant 0964507 and in part by the NetApp Faculty Fellowship program. The work of V. Lalitha is supported by a TCS Research Scholarship.

The authors of [7] handle multiple erasures by protecting each message symbol with a local code of length  $\leq r + \delta - 1$  and minimum distance  $\geq \delta$ , and derive the upper bound

$$d_{\min} \leq n - k + 1 - \left( \left\lceil \frac{k}{r} \right\rceil - 1 \right) (\delta - 1). \quad (2)$$

Pyramid codes, are once again shown to be optimal. This notion of locality is extended to the case when all code symbols are so protected and the existence of optimal codes with all-symbol locality is shown for the case when  $(r + \delta - 1)|n$ . Codes having the capability of locally recovering a failed node in the presence of any  $\delta - 1$  other node failures are also considered in [8]. Constructions based on partial geometry are provided and their rates computed.

A third approach to handling multiple erasures is presented in [9] in which the authors seek to protect each of the  $k$  message symbols by  $\delta - 1$  support-disjoint local parities, each of length  $\leq r + 1$ . The following upper bound is derived:

$$d_{\min} \leq n - k + 1 - \left( \left\lceil \frac{(k - 1)(\delta - 1) + 1}{(r - 1)(\delta - 1) + 1} \right\rceil - 1 \right), \quad (3)$$

and the existence of optimal codes is established for the case when  $n \geq k(r(\delta - 1) + 1)$ . The setting is extended to codes with all-symbol locality for handling 2 erasures, as well. A square-code construction that achieves the bound in (3) for restricted values of the code dimension  $k$  is presented. A related setting appears in [6], where once again  $\delta - 1$  support-disjoint local parities are used for the protection of all code symbols. Here however, the local parities are permitted to have different lengths. A key feature of this work is that the authors provide constructions of codes in which the code alphabet is small, on the order of the code length. Lower bounds to the minimum distance of the codes constructed are also provided.

A common underlying theme of the prior approaches in [7], [8], [9], [6] is that they implicitly assume the need for the recovery of multiple erased symbols in parallel. However, the need for locality does not preclude a sequential approach such as is adopted here. The sequential approach places a less-stringent requirement on the code and potentially allows us to construct codes with improved minimum distance while still enabling local recovery from erasures. In addition, the minimum distance bound derived here is universal in the sense that no code which permits all-symbol local recovery from 2 erasures can have larger distance regardless of approach adopted. The exact formulation and our approach to solving the problem are presented in Section III.

## B. Other Related Work

Explicit constructions of optimal codes with all-symbol locality for the single erasure case are provided in [10], [11], respectively based on Gabidullin maximum rank-distance and Reed-Solomon codes. Families of codes with all-symbol locality with small alphabet size (low field size) are constructed in [6]. Locality in the context of non-linear codes is considered in [4]. Codes with local regeneration are considered in [12], [13], [14]. Studies on the implementation and performance evaluation of codes with locality can be found in [5], [15].

Section II provides background on generalized Hamming weights (GHW). Our formulation and approach to the problem are outlined in Section III. An important connection between the  $k$ -cores of [2] and GHW is made in Section IV. The upper bound on  $d_{\min}$  and optimal code constructions can be found in Sections V and VI respectively. The final section, Section VII presents the analogous  $d_{\min}$  bound for the single-erasure case. The proofs of most statements appear in the Appendix.

## II. GENERALIZED HAMMING WEIGHTS

*Definition 1:* The  $i^{th}$ ,  $1 \leq i \leq k$ , GHW [16], [17] of an  $[n, k]$  code  $\mathcal{C}$  is the cardinality of the minimum support of an  $i$ -dimensional subcode of  $\mathcal{C}$ , i.e.,

$$d_i(\mathcal{C}) = d_i = \min_{\substack{\mathcal{D} < \mathcal{C} \\ \dim(\mathcal{D})=i}} |\text{supp}(\mathcal{D})|, \quad (4)$$

where  $\mathcal{D} < \mathcal{C}$ , is used to denote a subcode  $\mathcal{D}$  of  $\mathcal{C}$  and  $\text{supp}(\mathcal{D}) \triangleq \cup_{\mathbf{c} \in \mathcal{D}} \text{supp}(\mathbf{c})$ .

It is well known that

$$d_{\min}(\mathcal{C}) = d_1 < d_2 < \dots < d_k = n. \quad (5)$$

The complement of the set  $\{d_i, 1 \leq i \leq k\}$ , in  $[n]$ , will be termed as the set of *gap numbers* (more simply, gaps) of the code  $\mathcal{C}$  and denoted by  $\{g_i, 1 \leq i \leq n - k\}$ , i.e.,

$$\{g_i, 1 \leq i \leq n - k\} = [n] \setminus \{d_i, 1 \leq i \leq k\}. \quad (6)$$

Similarly, let  $\{d_j^\perp, 1 \leq j \leq n - k\}$  and  $\{g_i^\perp, 1 \leq i \leq k\}$  denote the GHWs and gaps of the dual code  $\mathcal{C}^\perp$ . The lemma below [16] relates the GHWs of  $\mathcal{C}$  to the gaps of  $\mathcal{C}^\perp$ .

*Lemma 2.1:*

$$\begin{aligned} d_i &= (n + 1) - g_{k-i+1}^\perp, \quad 1 \leq i \leq k, \\ d_{\min}(\mathcal{C}) &= d_1 = (n + 1) - g_k^\perp. \end{aligned}$$

We also note that if  $\mathcal{B}_0 < \mathcal{C}^\perp$ , then  $d_i(\mathcal{B}_0) \geq d_i(\mathcal{C}^\perp)$ , which implies that  $g_k^\perp \geq g_k(\mathcal{B}_0)$ . We thus obtain:

$$d_{\min}(\mathcal{C}) \leq n + 1 - g_k(\mathcal{B}_0). \quad (7)$$

### III. APPROACH AND RESULTS

Our focus in this paper, is on codes with all-symbol locality for the two-erasure case.

*Definition 2:* A code  $\mathcal{C}$  will be said to be *locally 2-reconstructible with locality  $r$* , if for any pair of code-symbol erasures, there exists a sequence of two local (and linear) parity-check computations that can be used to recover the erased symbols. By a local parity-check computation, we mean recovery via a parity whose support covers the coordinate being recovered and involves at most  $r$  other code symbols.

Note that under the above definition, it is permissible that the symbol recovered by the first local parity belong to the set of  $r$  symbols accessed by the second local parity. The families of all-symbol locality codes constructed in [7], [8], [9] for the case  $\delta = 3$  may all be regarded as examples of locally 2-reconstructible codes. In the sequel, we will refer to a locally 2-reconstructible code with locality  $r$  simply as a locally reconstructible code. Our principal results are an upper bound on the minimum distance of locally reconstructible codes and optimal constructions for a large class of code parameters. The steps involved in the derivation of the upper bound on  $d_{\min}$  are outlined below.

Given a locally reconstructible code  $\mathcal{C}$ , let  $\mathcal{B}_0$  denote the subcode of the dual code  $\mathcal{C}^\perp$ , spanned by all codewords  $\mathbf{c} \in \mathcal{C}^\perp$  of Hamming weight less than or equal to  $r + 1$ , i.e.,

$$\mathcal{B}_0 = \text{span} \left( \mathbf{c} \in \mathcal{C}^\perp, |\text{supp}(\mathbf{c})| \leq r + 1 \right). \quad (8)$$

*a) Step 1:* We first establish that the dimension  $b$  of  $\mathcal{B}_0$  satisfies the lower bound  $b \geq \frac{2n}{r+2}$ .

*b) Step 2:* Next, we observe from (7) that the minimum distance of the code  $\mathcal{C}$  satisfies  $d_{\min}(\mathcal{C}) \leq n + 1 - g_k(\mathcal{B}_0)$ .

*c) Step 3:* We then obtain a lower bound on the  $k^{\text{th}}$  gap of  $\mathcal{B}_0$  of the form  $g_k(\mathcal{B}_0) \geq \gamma_k$ , leading to the desired upper bound  $d_{\min}(\mathcal{C}) \leq n + 1 - \gamma_k$ . This step makes use of the lower bound on the dimension  $b$  of  $\mathcal{B}_0$ , derived in Step 1.

The same sequence of steps is also applied in Section VII to the case of codes with all-symbol locality for the single-erasure case. This results in a new bound on  $d_{\min}$  for this class of codes, tighter in general than that given by (1).

#### A. Optimal Constructions

We provide code constructions that are optimal with respect to the bound on  $d_{\min}$  given above in Step 3 whenever the block length  $n$  is of the form  $n = \frac{(r+\beta)(r+2)}{2}$ ,  $1 \leq \beta \leq r$ , with  $\beta|r$ . The steps involved are described below.

*a) Step 1:* We begin by constructing a code  $\mathcal{B}_0$  such that (i) the code formed by the null space of  $\mathcal{B}_0$  possesses the locally reconstructible property, (ii)  $\dim(\mathcal{B}_0) = b = \frac{2n}{r+2}$ , and (iii) the lower bound on the  $k^{\text{th}}$  gap is also achieved, i.e.,  $g_k(\mathcal{B}_0) = \gamma_k$ . Our construction of  $\mathcal{B}_0$  is based on Turán graphs [18], depends only on code length  $n$  and locality parameter  $r$ , and is independent of the dimension  $k$  of the desired code  $\mathcal{C}$ .

b) *Step 2*: Given the code  $\mathcal{B}_0$ , it turns out that it is always possible to find an  $[n, k]$  code  $\mathcal{C}$  such that  $\mathcal{B}_0$  is a subcode of  $\mathcal{C}^\perp$ ,  $g_k(\mathcal{C}^\perp) = g_k(\mathcal{B}_0)$  and this code  $\mathcal{C}$  is then the desired locally reconstructible code. It has the best possible minimum distance given by  $d_{\min}(\mathcal{C}) = n + 1 - \gamma_k$ .

This proof is an instance of a more general result that is important in its own right and which can potentially be applied in other situations as well. It combines the notion of a  $k$ -core introduced in [2] with the GHW structure of a code to enable construction of the best possible code of a given dimension when the code is linearly constrained. This is discussed in more detail in the next section.

#### IV. $k$ -CORES AND CONNECTION WITH GHWs

*Definition 3 ([2]):* Given a linear code  $\mathcal{B}_0$ , a set  $S \subseteq [n], |S| = \ell$  is termed an  $\ell$ -core of  $\mathcal{B}_0$  if  $\text{supp}(\mathbf{c}) \not\subseteq S$ ,  $\forall \mathbf{c} \in \mathcal{B}_0$ .

The above definition is equivalent to saying that  $\text{rank}(G'|_S) = \ell$ , for any  $S$  which is an  $\ell$ -core of  $\mathcal{B}_0$ , where  $G'$  denotes a generator matrix of  $\mathcal{B}_0^\perp$ . The lemma below was used in [2] to show the existence of all-symbol locality codes when  $(r+1)|n$ , and will also prove very useful here.

*Lemma 4.1 ([2]):* Let  $\mathcal{B}_0$  denote an  $[n, t]$  code over  $\mathbb{F}_q$ . Then for any  $k$  such that  $k \leq n - t$ , there exists an  $[n, k]$  code  $\mathcal{C}$  over  $\mathbb{F}_q$  such that

- (a)  $\mathcal{B}_0 < \mathcal{C}^\perp$ , and
  - (b) any  $S$  which is a  $k$ -core of  $\mathcal{B}_0$  is also a  $k$ -core of  $\mathcal{C}^\perp$ ,
- whenever  $q > kn^k$ .

In the following theorem, we obtain an expression for the minimum distance of the code whose existence is guaranteed by the above lemma.

*Theorem 4.2:* Let  $\mathcal{B}_0$  denote an  $[n, t]$  code and let  $\mathcal{C}$  denote an  $[n, k]$  code,  $k \leq n - t$ , such that

- (a)  $\mathcal{B}_0 < \mathcal{C}^\perp$ , and
- (b) any  $S$  which is a  $k$ -core of  $\mathcal{B}_0$  is also a  $k$ -core of  $\mathcal{C}^\perp$ .

The minimum distance of the code  $\mathcal{C}$  is given by

$$d_{\min}(\mathcal{C}) = n + 1 - g_k(\mathcal{B}_0). \quad (9)$$

*Proof:* See Appendix A. ■

Note from (7) that the minimum distance of the code  $\mathcal{C}$ , whenever  $\mathcal{B}_0 < \mathcal{C}^\perp$ , cannot be any larger than  $n + 1 - g_k(\mathcal{B}_0)$ . In addition to showing that the  $k^{\text{th}}$  gap of  $\mathcal{C}^\perp$  is same as that of  $\mathcal{B}_0$ , it is possible to identify all the GHWs of  $\mathcal{C}^\perp$  in terms of the GHWs of  $\mathcal{B}_0$ . This is stated in the following theorem.

*Theorem 4.3:* Let  $\mathcal{B}_0$  denote an  $[n, t]$  code and let  $\mathcal{C}$  denote an  $[n, k]$  code,  $k \leq n - t$ , such that

- (a)  $\mathcal{B}_0 < \mathcal{C}^\perp$ , and
- (b) any  $S$  which is a  $k$ -core of  $\mathcal{B}_0$  is also a  $k$ -core of  $\mathcal{C}^\perp$ .

The generalized Hamming weights of  $\mathcal{C}^\perp$  are given by

$$d_i(\mathcal{C}^\perp) = \begin{cases} d_i(\mathcal{B}_0), & 1 \leq i \leq g_k(\mathcal{B}_0) - k \\ i + k, & g_k(\mathcal{B}_0) - k + 1 \leq i \leq n - k \end{cases} \quad (10)$$

*Proof:* See Appendix B. ■

An illustration of (10) is given in Fig. 1 with parameters  $n = 15$ ,  $t = 5$  and  $k = 8$ . We see that the largest gap of  $\mathcal{C}^\perp$  is same as the  $k^{\text{th}}$  gap of  $\mathcal{B}_0$ . Moreover the GHWs of  $\mathcal{C}^\perp$  which appear to the left of the  $k^{\text{th}}$  gap are exactly same as those of  $\mathcal{B}_0$ .

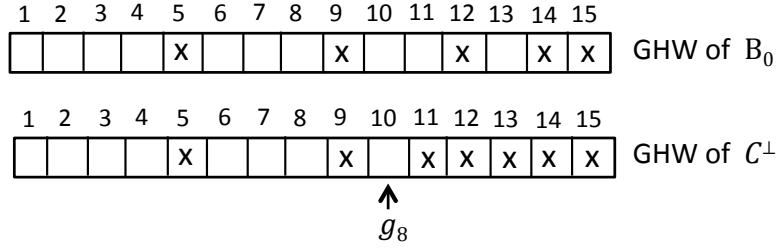


Fig. 1. Relation between GHWs of  $C^\perp$  and  $B_0$ , with  $n = 15$ ,  $t = 5$  and  $k = 8$ . GHWs are indicated by an 'X', gaps by a blank.

## V. MINIMUM DISTANCE BOUND FOR LOCALLY RECONSTRUCTIBLE CODES

In this section, we will obtain upper bounds on the GHWs of the subcode  $B_0$ , as defined in (8). This in turn will establish a lower bound on the  $k^{\text{th}}$  gap  $g_k(B_0)$  from which we will obtain an upper bound on the minimum distance of  $C$ . We begin with a characterization of a locally reconstructible code.

*Lemma 5.1:* Let  $\mathcal{A}_i$  denote the collection of all the local parities which cover the code symbol  $c_i$ ,  $1 \leq i \leq n$ . Code  $C$  is locally reconstructible iff (i)  $|\mathcal{A}_i| \geq 1$  and (ii)  $\mathcal{A}_i \neq \mathcal{A}_j$ ,  $\forall i, j$ ,  $i \neq j$ .

*Proof:* Straightforward. ■

The parallel recovery of two code symbols  $c_i$  and  $c_j$  is possible iff  $\mathcal{A}_i \not\subseteq \mathcal{A}_j$  and  $\mathcal{A}_j \not\subseteq \mathcal{A}_i$ . In the event that  $\mathcal{A}_i \subsetneq \mathcal{A}_j$ ,  $c_j$  can be recovered first through a local computation not involving  $c_i$  and having recovered  $c_j$ ,  $c_i$  can then be recovered.

*Theorem 5.2:* The dimension of the subcode  $B_0$  defined in (8) is lower bounded by

$$\dim(B_0) \geq \frac{2n}{r+2}. \quad (11)$$

*Proof:* See Appendix C. ■

*Corollary 5.3:* The rate of any code  $C$  which is locally reconstructible is upper bounded by

$$\frac{k}{n} \leq \frac{r}{r+2}. \quad (12)$$

The following lemma will be used to establish upper bounds on the GHWs of  $B_0$ .

*Lemma 5.4:* Let  $T$  be any set such that  $|T| = n \geq r+1$  and let  $S_i$ ,  $1 \leq i \leq b$  be subsets of  $T$  such that (i)  $\cup_{i=1}^b S_i = T$  and (ii)  $|S_i| = r+1$ ,  $\forall i \in \{1, 2, \dots, b\}$ . Define

$$f_m = \min_{\substack{I \subseteq [b] \\ |I|=m}} |\cup_{i \in I} S_i|, 1 \leq m \leq b. \quad (13)$$

Then,  $\forall m \in [b]$ ,  $f_m \leq e_m$ , where the  $\{e_m\}$  are obtained recursively as follows:

$$e_b = n, \quad (14)$$

$$e_{m-1} = e_m - \left\lceil \frac{2e_m}{m} \right\rceil + (r+1), \quad 2 \leq m \leq b. \quad (15)$$

*Proof:* See Appendix D. ■

Note that in Lemma 5.4, since  $\cup_{i=1}^b S_i = T$ , we have that  $b \geq \frac{n}{r+1}$  and thus setting  $m = b$  in (15) and dropping the ceiling function, we obtain

$$\begin{aligned} e_{b-1} &\leq \left(1 - \frac{2}{b}\right)e_b + (r+1) = \left(1 - \frac{2}{b}\right)n + (r+1), \\ &\leq \left(1 - \frac{2}{b}\right)b(r+1) + (r+1) = (b-1)(r+1). \end{aligned} \quad (16)$$

The arguments in (16) can be iterated further (with  $m = b - 1$  and so on) and from this it follows that Lemma 5.4 indeed implies the obvious bound  $f_m \leq m(r + 1)$ ,  $1 \leq m \leq b$ . In general the bounds given by Lemma 5.4 will be tighter than this and will take into account the fact that the total support has cardinality only  $n$ .

*Theorem 5.5:* Let  $\mathcal{C}$  denote an  $[n, k]$  locally reconstructible code and let  $\mathcal{B}_0$  be the subcode as defined in (8). Set  $b = \left\lceil \frac{2n}{r+2} \right\rceil$ . Then the first  $b$  GHWs of  $\mathcal{B}_0$ , and hence those of  $\mathcal{C}^\perp$ , are upper bounded by  $d_m(\mathcal{B}_0) \leq e_m$ ,  $1 \leq m \leq b$ , where  $e_m$  is as defined by (14) and (15). Furthermore, if  $\ell$  denotes the unique integer satisfying  $e_\ell < k + \ell < e_{\ell+1}$ , then the minimum distance of  $\mathcal{C}$  is upper bounded by

$$d_{\min}(\mathcal{C}) \leq n + 1 - (k + \ell). \quad (17)$$

*Proof:* Consider a basis of  $\mathcal{B}_0$  which are composed only of codewords of Hamming weight less than or equal to  $r + 1$ . Let  $\{S_i, i = 1, \dots, b\}$  denote their supports, where  $b \geq \frac{2n}{r+2}$ . The bounds on the GHWs of  $\mathcal{B}_0$  now follows directly upon applying Lemma 5.4 to the sets  $\{S_i, i = 1, \dots, b\}$ . (Note that if any set  $S_i$  has cardinality less than  $r + 1$ , one can simply substitute  $S_i$  with any set  $S'_i$  such that  $|S'_i| = r + 1$ ,  $S_i \subseteq S'_i$  and then apply Lemma 5.4.).

Given the bounds on the GHWs of  $\mathcal{B}_0$ , the  $k^{\text{th}}$  gap of  $\mathcal{B}_0$  is lower bounded by  $g_k(\mathcal{B}_0) \geq k + \ell$ , where  $\ell$  denotes the unique integer such that  $e_\ell < k + \ell < e_{\ell+1}$  (to see this, assume that first  $b$  GHWs are given exactly by the sequence  $\{e_m, m = 1, \dots, b\}$  and using this, identify the  $k^{\text{th}}$  gap). The bound on  $d_{\min}$  finally follows from (7). ■

A code  $\mathcal{C}$  will be called an *optimal locally reconstructible code* if it achieves the bound in (17) with equality.

## VI. OPTIMAL LOCALLY RECONSTRUCTIBLE CODES

In this section, we will describe a construction for optimal locally reconstructible codes for the case when the length of the code takes on the form

$$n = \frac{(r + \beta)(r + 2)}{2}, \quad (18)$$

with  $1 \leq \beta \leq r$  and  $\beta | r$ . The only restriction on the dimension  $k$  is the necessary rate restriction given by Corollary 5.3, i.e.,  $k \leq \frac{rn}{r+2}$ . As described in Section I, our approach to optimal code construction will involve first constructing a code  $\mathcal{B}_0$  which depends only on  $n, r$  and is independent of  $k$ . The construction of  $\mathcal{B}_0$  will be based on Turán graphs and will be such that

- (a)  $\mathcal{B}_0^\perp$  is locally reconstructible,
- (b)  $\dim(\mathcal{B}_0) = b = \frac{2n}{r+2} = r + \beta$ , and
- (c) all the  $b$  GHWs of  $\mathcal{B}_0$  achieve the upper bounds given by Theorem 5.5.

Once we have the code  $\mathcal{B}_0$ , the desired  $[n, k]$  code is simply the code  $\mathcal{C}$ , whose existence is guaranteed by Lemma 4.1. It is clear, based on the discussion in Section V and from Theorem 4.3 that this code  $\mathcal{C}$  will be an optimal locally reconstructible code.

### A. Construction of $\mathcal{B}_0$ Using Turán Graphs

Consider a graph with  $b = \frac{2n}{r+2} = r + \beta$  vertices. We partition the vertices into  $x = \frac{r+\beta}{\beta}$  partitions, each partition containing  $\beta$  vertices. We next place exactly one edge between any two vertices belonging to two distinct partitions. The resulting graph is known as a Turán graph on  $b$  vertices with  $x$  vertex partitions. The number of edges in this graph is  $\frac{x(x-1)\beta^2}{2} = n - b$  and each vertex is connected to exactly  $(x - 1)\beta = r$  other vertices. Let the vertices be labelled from 1 to  $b$  and the edges be labelled from  $b + 1$  to  $n$ , without paying attention to order.

To convert the graph into a code, we proceed as follows. Associate a local parity with each of the  $b$  vertices, let parity  $\underline{p}_i$  be associated with vertex  $i$ ,  $1 \leq i \leq b$ . Let  $\{i_1, i_2, \dots, i_r\}$  denote all the edges which are incident up on vertex  $i$ . Then, the support  $S_i \subseteq [n]$  of the local parity  $\underline{p}_i$  is set as

$$S_i = \{i, i_1, i_2, \dots, i_r\} \quad (19)$$

and the codeword  $\mathbf{c}_i$  corresponding to  $\underline{p}_i$  is identified as the all-1 vector in these  $r + 1$  coordinates (with zeros in the remaining  $n - (r + 1)$  coordinates). Set  $\mathcal{B}_0 = \text{span}(\mathbf{c}_i, 1 \leq i \leq b)$ . It is easily verified that the code  $\mathcal{B}_0^\perp$  is locally reconstructible and that its dual  $\mathcal{B}_0$  has dimension  $b = \frac{2n}{r+2} = r + \beta$ . Before proceeding to evaluate the

GHWs of the code  $\mathcal{B}_0$  and proving their optimality w.r.t. Theorem 5.5, we first illustrate the construction using two examples.

*Example 1:* Consider the parameters  $r = 3$  and  $\beta = 1$ , which implies that the length  $n = 10$ . When  $\beta = 1$ , note that the number of partitions  $x = r + 1 = 4$ , and each partition has just one vertex. Thus the total number of vertices  $b = r + 1 = 4$  and the graph is simply a completely connected graph on  $b = 4$  vertices. The generator matrix of the code  $\mathcal{B}_0$  (upto permutation of columns) in this case is given by

$$H_0 = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & & & 1 & 1 & 1 \\ & 1 & & 1 & & 1 & 1 \\ & & 1 & & 1 & 1 & 1 \end{bmatrix}.$$

Note that the sum of the four rows of  $H_0$  gives a local parity having support  $\{4, 7, 9, 10\}$ . Thus we see that the code  $\mathcal{B}_0$  guarantees that any code symbol is covered by two support disjoint local parities. In general, when  $\beta = 1$ , the construction gives  $(r, \delta = 3)_c$  all-symbol locality codes, with length  $n = \frac{(r+1)(r+2)}{2}$  (see Section I for a definition of  $(r, \delta)_c$  codes).

*Example 2:* In this example, let  $\beta = r = 3$ , which implies that the length  $n = 15$ . When  $\beta = r$ , we get a bipartite graph with  $r$  vertices on each of the two partitions. The generator matrix of the code  $\mathcal{B}_0$ , which is the span of 6 codewords, is given (up to permutation of columns) by

$$H_0 = \begin{bmatrix} 1 & 1 & 1 & 1 & & & & & & & & & & & \\ & & & & 1 & 1 & 1 & 1 & & & & & & & \\ 1 & & & & & & & & 1 & 1 & 1 & 1 & & & \\ & 1 & & & & & & & 1 & & & & & 1 & \\ & & 1 & & & 1 & & & & 1 & & & & 1 & \\ & & & 1 & & & 1 & & & & 1 & & & & 1 \end{bmatrix}.$$

It can be verified that any non-trivial linear combination (resulting in vectors other than those appearing in the rows of  $H_0$ ) of the 6 vectors results in a codeword whose Hamming weight  $\geq 5$ . As a result, each of the code symbols  $\{c_4, c_8, c_{12}, c_{13}, c_{14}, c_{15}\}$  is covered by only one local parity and thus parallel decoding of two erased symbols may not always be possible. For instance, if symbols  $c_3, c_4$  get erased, we must necessarily decode  $c_3$  first before decoding  $c_4$ .

The Turán graph construction, when  $\beta = r$  is closely related to the square code construction presented in [9]. The square code construction was used to guarantee two support disjoint parities for each code word symbol. For the current example, the closest relative from the square code family has length 16 and has one more local parity (in addition to those described by  $\mathcal{B}_0$ ) covering the coordinates  $\{4, 8, 12, 16\}$ . A second local parity which covers  $c_{16}$  can be obtained as a linear combination of all the 7 parities, and this will be a parity on the support  $\{13, 14, 15, 16\}$ .

## B. Generalized Hamming Weights of the Constructed Code $\mathcal{B}_0$

*Theorem 6.1:* Consider the code  $\mathcal{B}_0$  obtained via the Turán graph construction along with support sets  $\{S_i, 1 \leq i \leq b = r + \beta\}$ , as described in (19), associated to the  $b$  local parities  $\{p_i\}$ . Then the sets  $\{S_i, 1 \leq i \leq b = r + \beta\}$  achieve the upper bounds given in Lemma 5.4 with equality, i.e.,  $\forall m \in [b], f_m = e_m$ , where  $f_m$  is as described by (13) and  $e_m$  is as defined recursively by (14) and (15).

*Proof:* See Appendix E. ■

We use the following lemma to argue that the  $m^{\text{th}}$  GHW of  $\mathcal{B}_0$  is indeed given by  $f_m$  i.e., any other  $m$  dimensional subspace of  $\mathcal{B}_0$  (i.e., other than those generated by  $m$  subsets of the basis vectors) will have a support whose cardinality is no less than  $f_m$ .

*Lemma 6.2:* Let  $\mathcal{D}$  denote an  $[n, t]$  linear code and let  $\{\mathbf{v}_1, \dots, \mathbf{v}_t\}$  be a basis for the code  $\mathcal{D}$ . Also, let  $R_i = \text{supp}(\mathbf{v}_i)$  and suppose that the sets  $\{R_i\}$  are such that

- (a)  $|R_i \cap R_j| \leq 1, \forall i, j, i \neq j$ ,
- (b) any element  $\ell \in [n]$  belongs to at most two sets among the sets  $\{R_i\}$ , and
- (c)  $|R_i \setminus \bigcup_{\substack{j=1 \\ j \neq i}}^t R_j| \geq 1$ .

Then the generalized Hamming weights of the code  $\mathcal{D}$  are given by

$$d_m(\mathcal{D}) = \min_{\substack{\mathcal{I} \subseteq [t] \\ |\mathcal{I}|=m}} |\cup_{i \in \mathcal{I}} R_i|, \quad 1 \leq m \leq t. \quad (20)$$

*Proof:* See Appendix F. ■

We now note that Lemma 6.2 is readily applicable to the code  $\mathcal{B}_0$  obtained via the Turán graph construction. From this we conclude that the GHWs of this code  $\mathcal{B}_0$  achieve the upper bounds given by Theorem 5.5.

## VII. A NEW UPPER BOUND ON MINIMUM DISTANCE FOR THE SINGLE ERASURE CASE

The approach described in Section V directly applies to the setting of codes with all-symbol locality which can handle single erasures. This results in a new upper bound on  $d_{\min}$  for this class of codes which is in general tighter than that given by (1). Let  $\mathcal{C}$  be an  $[n, k, d_{\min}]$  code having  $(r, \delta = 2)$  all-symbol locality, i.e., any code symbol is covered by a local parity. As with locally reconstructible codes, consider the subcode  $\mathcal{B}_0$  of  $\mathcal{C}^\perp$  which is obtained as the span of all codewords of Hamming weight less than or equal to  $r+1$ , i.e.,  $\mathcal{B}_0 = \text{span}(\mathbf{c} \in \mathcal{C}^\perp, |\text{supp}(\mathbf{c})| \leq r+1)$ . It is easy to see that  $\dim(\mathcal{B}_0) \geq \frac{n}{r+1}$ . Lemma 5.4 can now be applied to this  $\mathcal{B}_0$  which enables us to upper bound the GHWs of  $\mathcal{B}_0$  and in turn, upper bound the minimum distance of  $\mathcal{C}$ .

*Theorem 7.1:* Let  $b = \left\lceil \frac{n}{r+1} \right\rceil$ . Then, the first  $b$  generalized Hamming weights of the subcode  $\mathcal{B}_0$  defined above are upper bounded by  $d_m(\mathcal{B}_0) \leq e_m$ ,  $1 \leq m \leq b$ , where  $e_m$  is as recursively defined by  $e_b = n$ , and

$$e_{m-1} = e_m - \left\lceil \frac{2e_m}{m} \right\rceil + (r+1), \quad 2 \leq m \leq b. \quad (21)$$

Furthermore, if we let  $\ell$  to denote the unique integer such that  $e_\ell < k + \ell < e_{\ell+1}$ , the minimum distance of the all-symbol locality code  $\mathcal{C}$  is upper bounded by

$$d_{\min}(\mathcal{C}) \leq n + 1 - (k + \ell). \quad (22)$$

*Proof:* Similar to the proof of Theorem 5.5. ■

In order to compare the upper bound given by (22) with that given by (1), we note the bound given by (1) can be obtained by first upper bounding the GHWs of  $\mathcal{B}_0$  by

$$d_m(\mathcal{B}_0) \leq m(r+1), \quad 1 \leq m \leq b-1, \quad \text{and } d_b(\mathcal{B}_0) \leq n, \quad (23)$$

where  $b = \left\lceil \frac{n}{r+1} \right\rceil$ , and then calculating the  $k^{\text{th}}$  gap based on these bounds. But from the discussion in Section V (see (16)), we know that the bounds on GHWs of  $\mathcal{B}_0$  given by Theorem 7.1 are, in general, tighter than the bounds in (23) and hence we conclude that the minimum distance bound given by (22) is also tighter, in general, than that given by (1). We would, however, like to remark that it is always possible [6] to achieve a minimum distance which is at most one less than that suggested by the upper bound in (1). In Fig. 2, we plot the two bounds as a function of dimension  $k$ , for the case when  $n = 18$  and  $r = 3$ .

## REFERENCES

- [1] A. G. Dimakis, P. B. Godfrey, Y. Wu, M. J. Wainwright, and K. Ramchandran, "Network coding for distributed storage systems," *IEEE Trans. Inf. Theory*, vol. 56, no. 9, pp. 4539–4551, Sep. 2010.
- [2] P. Gopalan, C. Huang, H. Simitci, and S. Yekhanin, "On the Locality of Codeword Symbols," *IEEE Trans. Inf. Theory*, vol. 58, no. 11, pp. 6925–6934, Nov. 2012.
- [3] C. Huang, M. Chen, and J. Li, "Pyramid codes: Flexible schemes to trade space for access efficiency in reliable data storage systems," in *Proc. 6th IEEE Int. Symposium on Network Computing and Applications (NCA)*, 2007, pp. 79–86.
- [4] D. S. Papailiopoulos and A. G. Dimakis, "Locally repairable codes," in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, Cambridge, MA, Jul. 2012, pp. 2771–2775.
- [5] C. Huang, H. Simitci, Y. Xu, A. Ogus, B. Calder, P. Gopalan, J. Li, and S. Yekhanin, "Erasure coding in windows azure storage," in *Proc. 2012 USENIX Annual Technical Conference (ATC)*, Boston, MA, 2012, pp. 15–26.
- [6] I. Tamo and A. Barg, "A family of optimal locally recoverable codes," 2013. [Online]. Available: arXiv:1311.3284



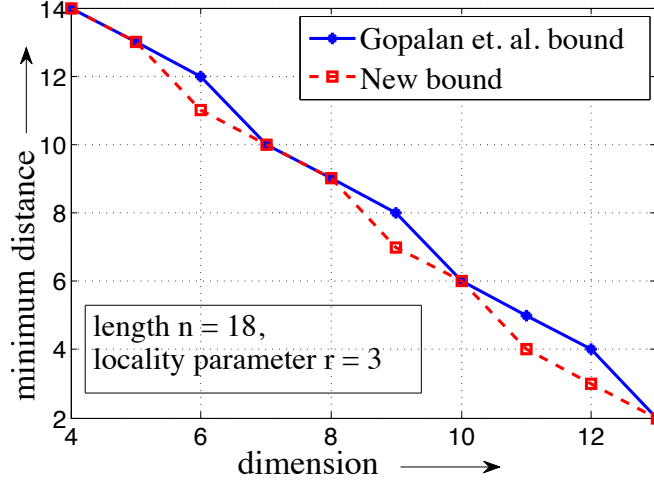


Fig. 2. Comparing the bounds on  $d_{\min}$  for varying  $k$  for codes with information and all-symbol locality, with  $n = 18$  and  $r = 3$ .

- [7] N. Prakash, G. M. Kamath, V. Lalitha, and P. V. Kumar, "Optimal linear codes with a local-error-correction property," in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, Cambridge, MA, Jul. 2012, pp. 2776–2780.
- [8] L. Parnes-Juarez, H. D. L. Hollmann, and F. Oggier, "Locally repairable codes with multiple repair alternatives," in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, 2013, pp. 892–896.
- [9] A. Wang and Z. Zhang, "Repair locality with multiple erasure tolerance," *CoRR*, vol. abs/1306.4774, 2013.
- [10] N. Silberstein, A. S. Rawat, and S. Vishwanath, "Error resilience in distributed storage via rank-metric codes," in *Proc. 50th Annual Allerton Conf. on Communication, Control, and Computing (Allerton)*, Urbana-Champaign, IL, Oct. 2012, pp. 1150–1157.
- [11] I. Tamo, D. S. Papailiopoulos, and A. G. Dimakis, "Optimal locally repairable codes and connections to matroid theory," in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, 2013, pp. 1814–1818.
- [12] G. M. Kamath, N. Prakash, V. Lalitha, and P. V. Kumar, "Codes with local regeneration," in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, 2013, pp. 1606–1610.
- [13] N. Silberstein, A. S. Rawat, O. O. Koyluoglu, and S. Vishwanath, "Optimal locally repairable codes via rank-metric codes," in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, 2013, pp. 1819–1823.
- [14] G. M. Kamath, N. Silberstein, N. Prakash, A. S. Rawat, V. Lalitha, O. O. Koyluoglu, P. V. Kumar, and S. Vishwanath, "Explicit mbr all-symbol locality codes," in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, 2013, pp. 504–508.
- [15] M. Sathiamoorthy, M. Asteris, D. Papailiopoulos, A. G. Dimakis, R. Vadali, S. Chen, and D. Borthakur, "Xoring elephants: Novel erasure codes for big data," 2013. [Online]. Available: arXiv:1301.3791
- [16] V. K. Wei, "Generalized Hamming weights for linear codes," *IEEE Trans. Inf. Theory*, vol. 37, no. 5, pp. 1412–1418, 1991.
- [17] T. Helleseth, T. Klove, V.I. Levenshtein, and O. Ytrehus, "Bounds on the minimum support weights," *IEEE Trans. Inf. Theory*, vol. 41, no. 2, pp. 432–440, 1995.
- [18] R. Diestel, *Graph Theory*. Springer Verlag, 2005.

## APPENDIX A PROOF OF THEOREM 4.2

Let  $\mathcal{B}_0$  denote an  $[n, t]$  code and let  $\mathcal{C}$  denote an  $[n, k]$  code,  $k \leq n - t$ , such that

- (a)  $\mathcal{B}_0 < \mathcal{C}^\perp$ , and
- (b) any  $S$  which is a  $k$ -core of  $\mathcal{B}_0$  is also a  $k$ -core of  $\mathcal{C}^\perp$ .

We claim that for any set  $S \subseteq [n]$ ,  $|S| = g_k(\mathcal{B}_0)$ , there exists  $S' \subseteq S$ ,  $|S'| = k$  such that  $S'$  is a  $k$ -core of  $\mathcal{B}_0$ . Supposing that this is true, it then means that for any set  $S \subseteq [n]$ ,  $|S| = g_k(\mathcal{B}_0)$ ,  $\text{rank}(G|_S) = k$ . Clearly, this would imply that  $d_{\min}(\mathcal{C}) \geq n - |S| + 1 = n - g_k(\mathcal{B}_0) + 1$ . However, since  $d_{\min}(\mathcal{C}) = n - g_k(\mathcal{C}^\perp) + 1 \leq n - g_k(\mathcal{B}_0) + 1$ , we conclude that  $d_{\min}(\mathcal{C}) = n + 1 - g_k(\mathcal{B}_0)$ .

We will now prove our claim that for any set  $S \subseteq [n]$ ,  $|S| = g_k$ , there exists  $S' \subseteq S$ ,  $|S'| = k$  such that  $S'$  is a  $k$ -core of  $\mathcal{B}_0$ . Toward this, let  $\mathcal{C}'$  denote the code  $\mathcal{B}_0$  shortened to the set  $S$ . We note that  $\dim(\mathcal{C}') \leq g_k - k$ . To see why this is true, if we suppose that  $\dim(\mathcal{C}') > g_k - k$ , then it will imply that  $d_{g_k-k+1}(\mathcal{B}_0) \leq g_k$ , where  $d_{g_k-k+1}(\mathcal{B}_0)$  denotes the  $(g_k - k + 1)^{\text{th}}$  Generalized Hamming weight of  $\mathcal{B}_0$ . But this contradicts the fact that the number of GHWs of  $\mathcal{B}_0$  till  $g_k$  is exactly  $g_k - k$  and hence we get that  $\dim(\mathcal{C}') \leq g_k - k$ . Now, let  $\rho$  denote the dimension of  $\mathcal{C}'$  and let  $H' = [I_\rho | P_{\rho \times (g_k - \rho)}]$  denote a generator matrix of  $\mathcal{C}'$ , up to permutation of columns. If  $T$  denotes the support of the matrix  $P_{\rho \times (g_k - \rho)}$ , then any  $S' \subseteq T$ ,  $|S'| = k$  is  $k$ -core of  $\mathcal{B}_0$ .

APPENDIX B  
PROOF OF THEOREM 4.3

Let  $\mathcal{B}_0$  denote an  $[n, t]$  code and let  $\mathcal{C}$  denote an  $[n, k]$  code,  $k \leq n - t$ , such that

- (a)  $\mathcal{B}_0 < \mathcal{C}^\perp$ , and
- (b) any  $S$  which is a  $k$ -core of  $\mathcal{B}_0$  is also a  $k$ -core of  $\mathcal{C}^\perp$ .

Note that Theorem 4.2 implies that  $g_k(\mathcal{C}^\perp) = g_k(\mathcal{B}_0)$ . This determines the last  $n - g_k(\mathcal{B}_0)$  GHWs of  $\mathcal{C}^\perp$  and are given by

$$d_i(\mathcal{C}^\perp) = i + k, \quad g_k(\mathcal{B}_0) - k + 1 \leq i \leq n - k. \quad (24)$$

Assuming that  $g_k(\mathcal{B}_0) - k \geq 1$ , it now remains to be proved that

$$d_i(\mathcal{C}^\perp) = d_i(\mathcal{B}_0), \quad 1 \leq i \leq g_k(\mathcal{B}_0) - k, \quad (25)$$

i.e., the first  $g_k(\mathcal{B}_0) - k$  GHWs of  $\mathcal{C}^\perp$  are exactly same as those of  $\mathcal{B}_0$ . Toward this, we first note that any set  $S$  which is a  $b$ -core of  $\mathcal{B}_0$  is also a  $b$ -core of  $\mathcal{C}^\perp$ , for any  $b$  such that  $b < k$ . This is because for any  $S$  which is a  $b$ -core of  $\mathcal{B}_0$  ( $b < k$ ), there exists a  $k$ -core  $S'$  of  $\mathcal{B}_0$  such that  $S \subseteq S'$ . We also claim that for any set  $S \subseteq [n]$ ,  $|S| = g_b$ , there exists  $S' \subseteq S$ ,  $|S'| = b$  such that  $S'$  is a  $b$ -core of  $\mathcal{B}_0$ . Proof is similar to the claim regarding  $k$ -cores which appeared in the proof of Theorem 4.2.

We will now prove (25) via induction starting at  $i = 1$ . Let us denote  $\ell = g_k(\mathcal{B}_0) - k$ . Note that

$$k > (d_1(\mathcal{B}_0) - 1) + (d_2(\mathcal{B}_0) - d_1(\mathcal{B}_0) - 1) + \dots + (d_\ell(\mathcal{B}_0) - d_{\ell-1}(\mathcal{B}_0) - 1). \quad (26)$$

Now, suppose that  $d_1(\mathcal{C}^\perp) \leq d_1(\mathcal{B}_0) - 1$ . Clearly any  $S$  such that  $|S| = d_1(\mathcal{B}_0) - 1$  is an  $|S|$ -core of  $\mathcal{B}_0$ . From (26), we see that  $d_1(\mathcal{B}_0) - 1 < k$  and hence  $S$  is also an  $|S|$ -core of  $\mathcal{C}^\perp$ . But this contradicts the assumption that  $d_1(\mathcal{C}^\perp) \leq d_1(\mathcal{B}_0) - 1$  and hence we conclude that  $d_1(\mathcal{C}^\perp) = d_1(\mathcal{B}_0)$ . Next, assume that  $d_i(\mathcal{C}^\perp) = d_i(\mathcal{B}_0)$ , for some  $i$  such that  $1 \leq i \leq \ell - 1$ . We will now prove that  $d_{i+1}(\mathcal{C}^\perp) = d_{i+1}(\mathcal{B}_0)$ . We consider the following cases:

- (a)  $d_{i+1}(\mathcal{B}_0) = d_i(\mathcal{B}_0) + 1$ . In this case, note that

$$d_i(\mathcal{B}_0) \stackrel{(i)}{=} d_i(\mathcal{C}^\perp) \quad (27)$$

$$\stackrel{(ii)}{<} d_{i+1}(\mathcal{C}^\perp) \quad (28)$$

$$\stackrel{(iii)}{\leq} d_{i+1}(\mathcal{B}_0) \quad (29)$$

$$= d_i(\mathcal{B}_0) + 1, \quad (30)$$

where (i) follows from the induction hypothesis and (ii) follows from (5). This then implies that (iii) must be an equality, i.e.,  $d_{i+1}(\mathcal{C}^\perp) = d_{i+1}(\mathcal{B}_0)$ .

- (b)  $d_{i+1}(\mathcal{B}_0) > d_i(\mathcal{B}_0) + 1$ . In this case, if we let  $m = d_{i+1}(\mathcal{B}_0) - (i + 1)$ , note that

$$g_m(\mathcal{B}_0) = d_{i+1}(\mathcal{B}_0) - 1. \quad (31)$$

Now, if  $S$  is any set such that  $|S| = g_m(\mathcal{B}_0)$ , then there exists  $S' \subseteq S$ ,  $|S'| = m$  and  $S'$  is an  $m$ -core of  $\mathcal{B}_0$ . From (26), we see that  $m < k$  and hence  $S'$  is also an  $m$ -core of  $\mathcal{C}^\perp$ . Now, without loss of generality, suppose that  $d_{i+1}(\mathcal{C}^\perp) = d_{i+1}(\mathcal{B}_0) - 1$ . Also, let  $\mathcal{D}$  denote an  $(i + 1)$ -dimensional subcode of  $\mathcal{C}^\perp$  having support  $S_{\mathcal{D}}$  such that  $|S_{\mathcal{D}}| = d_{i+1}(\mathcal{C}^\perp)$ . Note that for any set  $T \subseteq S_{\mathcal{D}}$  such that  $|T| = |S_{\mathcal{D}}| - i = m$ , one can find a non-zero vector in  $\mathcal{D}$  whose support is fully contained within  $T$  and hence there cannot exist an  $m$ -core of  $\mathcal{C}^\perp$  within  $S_{\mathcal{D}}$ . However, we know that this is not true and hence we conclude that  $d_{i+1}(\mathcal{C}^\perp) = d_{i+1}(\mathcal{B}_0)$ .

APPENDIX C  
PROOF OF THEOREM 5.2

Consider the code

$$\mathcal{B}_0 = \text{span} \left( \mathbf{c} \in \mathcal{C}^\perp, |\text{supp}(\mathbf{c})| \leq r + 1 \right) \quad (32)$$

and let  $\mathcal{B} = \{\mathbf{c}_1, \dots, \mathbf{c}_b\}$  denote a basis for  $\mathcal{B}_0$  such that  $|\text{supp}(\mathbf{c}_i)| \leq r+1, \forall i \in [b]$ . Also, let  $S_i = \text{supp}(\mathbf{c}_i)$ . Define the quantity

$$s_i = \left| S_i \setminus \bigcup_{\substack{j=1 \\ j \neq i}}^b S_j \right|, \quad 1 \leq i \leq b. \quad (33)$$

We claim that for the code  $\mathcal{C}$  to be locally reconstructible, it must necessarily be true that  $s_i \leq 1, \forall i \in [b]$ . To see this, suppose that for some  $i$ ,  $s_i \geq 2$  and let  $\{\ell_1, \ell_2\} \subseteq S_i \setminus \bigcup_{\substack{j=1 \\ j \neq i}}^b S_j$ . Then, if  $A_{\ell_1}$  and  $A_{\ell_2}$ , respectively denote all the local parities covering the code symbols  $c_{\ell_1}$  and  $c_{\ell_2}$ , it would mean that  $A_{\ell_1} = A_{\ell_2}$ . This is because  $\mathcal{B}$  is a basis and any linear combination whose support contains  $\ell_1$  will also contain  $\ell_2$ . The claim now follows by noting from Lemma 5.1 that the code cannot be locally reconstructible unless  $A_{\ell_1} \neq A_{\ell_2}$ .

In order to proceed and complete the proof of the theorem, we note that  $n - \sum_{i=1}^b s_i$  code symbols are covered by more than one of the sets  $S_i, i \in [b]$  and hence it must be true that

$$\sum_{i=1}^b s_i + 2 \left( n - \sum_{i=1}^b s_i \right) \leq b(r+1) \quad (34)$$

$$\implies 2n - \sum_{i=1}^b s_i \leq b(r+1) \quad (35)$$

$$\implies 2n - b \leq b(r+1) \quad (36)$$

$$\implies b \geq \frac{2n}{r+2}, \quad (37)$$

where (36) follows since  $s_i \leq 1, \forall i \in [b]$ .

#### APPENDIX D PROOF OF LEMMA 5.4

We will prove the lemma via induction, starting at  $m = b$  and decrementing  $m$  at each step. Clearly  $f_b \leq e_b = n$ . Now assuming that for some  $m, 2 \leq m \leq b, f_m \leq e_m$ , we will prove that  $f_{m-1} \leq e_{m-1}$ . Without loss of generality, let  $\{S_i, 1 \leq i \leq m\}$  be such that  $|\bigcup_{i=1}^m S_i| = f_m$ . Define

$$s_i = \left| S_i \setminus \bigcup_{\substack{j=1 \\ j \neq i}}^m S_j \right|, \quad 1 \leq i \leq m. \quad (38)$$

Then, noting that  $n - \sum_{i=1}^m s_i$  elements are covered by more than one of the sets  $S_i, i \in [m]$ , we get that

$$\sum_{i=1}^m s_i + 2 \left( f_m - \sum_{i=1}^m s_i \right) \leq m(r+1) \quad (39)$$

$$\implies \sum_{i=1}^m s_i \geq 2f_m - m(r+1). \quad (40)$$

Also, let  $s^* = \max_{i \in [m]} s_i$  and without loss of generality, assume that  $s_1 = s^*$ . Note that in this case,  $s_1 \geq \frac{\sum_{i=1}^m s_i}{m}$ .

Now, if we consider union of the sets  $\{S_i, 2 \leq i \leq m\}$ , we get that

$$|\cup_{i=2}^m S_i| = f_m - s_1 \quad (41)$$

$$\leq f_m - \frac{\sum_{i=1}^m s_i}{m} \quad (42)$$

$$\leq f_m - \frac{2f_m - m(r+1)}{m} \quad (43)$$

$$= \frac{m-2}{m} f_m + (r+1) \quad (44)$$

$$\leq \frac{m-2}{m} e_m + (r+1) \quad (45)$$

$$= e_m - \frac{2e_m}{m} + (r+1), \quad (46)$$

$$(47)$$

which implies that  $f_{m-1} \leq e_m - \frac{2e_m}{m} + (r+1)$ . Finally, noting that  $f_{m-1}$  is an integer, we get that

$$f_{m-1} \leq e_m - \left\lceil \frac{2e_m}{m} \right\rceil + (r+1) = e_{m-1} \quad (48)$$

#### APPENDIX E PROOF OF THEOREM 6.1

Consider any set of  $m$  parities, say  $\{p_1, \dots, p_m\}$  having supports  $S_1, \dots, S_m$ . Note that by our construction the parity  $p_i$  corresponds to the vertex  $i$  in the Turán graph. The cardinality of union of the supports  $S_1, \dots, S_m$  can be calculated from the graph as  $|\cup_{i=1}^m S_i| = m + |E|$ , where  $E$  is the set of all the edges in the graph with at least one of the end points being a vertex belonging to the set  $[m]$ . The quantity  $|E|$  can be equivalently be computed by first counting the number of edges in the graph restricted to the remaining vertices  $\{m+1, m+2, \dots, r+\beta\}$  and then subtracting it from the total number of edges in the original graph. Thus, for calculating  $f_m$ , it is sufficient to find a restricted graph on  $r+\beta-m$  vertices having the maximum number of edges. Let  $r+\beta-m = ux+v, 0 \leq u \leq \beta-1, 0 \leq v \leq x-1$ . Then, it is easy to see that the number of edges in a restricted graph on  $r+\beta-m$  vertices is maximized if the restricted graph consists of  $u+1$  vertices each from any  $v$  out of the  $x$  partitions and  $u$  vertices each from the remaining  $x-v$  partitions.

It is straightforward to see now that the difference  $f_m - f_{m-1}, 2 \leq m \leq r+\beta$ , is given by

$$\begin{aligned} f_m - f_{m-1} &= (u+1)v + u(x-v-1) + 1 \\ &= v + ux - u + 1. \end{aligned} \quad (49)$$

The expression in (49) is evaluated for  $m, 2 \leq m \leq r+\beta$  and is shown in Fig. 3. For the array given in Fig. 3, we number the rows from 0 to  $\beta-1$  and the columns from 0 to  $x-1$ . Then, the value of  $f_m - f_{m-1}$  is simply the  $(u, v)^{\text{th}}$  entry in this array.

Now, we will show that the sequence  $\{f_m\}$  as defined by (49), satisfies the recursion given in (14) and (15), i.e.,  $f_{r+\beta} = n$  and

$$f_m - f_{m-1} = \left\lceil \frac{2f_m}{m} \right\rceil - (r+1). \quad (50)$$

Since the sequence  $\{e_m\}, 1 \leq m \leq b = r+\beta$  is unique (given that  $e_m = n$ ), it then follows that  $f_m = e_m, 1 \leq m \leq b = r+\beta$ , which will complete our proof. We begin by noting that

$$f_m = n - \sum_{i=r+\beta}^{m+1} (f_i - f_{i-1}). \quad (51)$$

$\beta$ rows	1	2	$\dots$		$x$
	$x$	$x+1$	$\dots$		$2x-1$
	$\vdots$				$\vdots$
					$(\beta-1)x-$ $(\beta-2)$
	$(\beta-1)x-$ $(\beta-2)$	$(\beta-1)x-$ $(\beta-1)$	$\dots$ $\dots$	$\beta(x-1)$	
$x$ columns					

Fig. 3.  $f_m - f_{m-1}$ ,  $r + \beta \geq m \geq 2$  of  $\mathcal{B}_0$  obtained via Turán graph construction, where the sequence of differences are in the descending order and the matrix has to be read row after row from left to right.

Next, note that the sum  $\sum_{i=r+\beta}^{m+1} (f_i - f_{i-1})$  can be calculated from the array in Fig. 3 as sum of the first  $r + \beta - m$  entries, where the entries are read from left to right in each row and the rows are read from top to bottom, i.e.,

$$\sum_{i=r+\beta}^{m+1} (f_i - f_{i-1}) = \sum_{i=1}^{ux-u+v} i + \sum_{i=1}^u (ix - i + 1), \quad (52)$$

where the first term on the R.H.S counts all the unique elements once and the second term on the R.H.S counts the repeated terms. Combining (51) and (52), we get that

$$\frac{2f_m}{m} = \frac{(r+\beta)(r+2) - (ux+v-u)(ux+v-u+1) - u(u+1)x + u(u-1)}{(r+\beta) - (ux+v)} \quad (53)$$

$$= \frac{[(r+\beta) - (ux+v)][(r+1) + (ux+v-u)] + (r+\beta) - ux - v\beta + uv}{(r+\beta) - (ux+v)} \quad (54)$$

$$= \frac{[(r+\beta) - (ux+v)]q' + r'}{(r+\beta) - (ux+v)}, \quad (55)$$

where  $q' = (r+1) + (ux+v-u)$  and  $r' = (r+\beta) - ux - v\beta + uv$ . It is straightforward to check that  $1 \leq r' \leq (r+\beta) - (ux+v)$ , which implies that

$$\left\lceil \frac{2f_m}{m} \right\rceil = q' + 1 = r + 1 + ux + v - u + 1. \quad (56)$$

Combining (49) with (56), we finally get that

$$f_m - f_{m-1} = \left\lceil \frac{2f_m}{m} \right\rceil - (r+1). \quad (57)$$

## APPENDIX F

### PROOF OF LEMMA 6.2

Consider an  $m$  dimensional subcode  $\mathcal{D}'$  of  $\mathcal{D}$  and let  $\mathcal{D}'$  have a basis  $\{\mathbf{u}_1, \dots, \mathbf{u}_m\}$ . Without loss of generality, let us assume that the basis of  $\mathcal{D}'$  is obtained as

$$\begin{bmatrix} \mathbf{u}_1 \\ \vdots \\ \mathbf{u}_m \end{bmatrix} = \begin{bmatrix} I_m & | & B_{m \times (b-m)} \end{bmatrix} \begin{bmatrix} \mathbf{v}_1 \\ \vdots \\ \mathbf{v}_m \\ \mathbf{v}_{m+1} \\ \vdots \\ \mathbf{v}_b \end{bmatrix}. \quad (58)$$

Also, let  $\{R'_i, 1 \leq i \leq m\}$  denote the supports of the vectors  $\{\mathbf{u}_1, \dots, \mathbf{u}_m\}$ . We claim that  $|\cup_{i=1}^m R'_i| \geq |\cup_{i=1}^m R_i|$ . To see this, we consider any element  $x \in \cup_{i=1}^m R_i$  and examine what happens to it when the  $m$  linear combinations are taken. We divide the discussion into the following cases:

- (a)  $x \in R_i, i \leq m$  and does not belong to any other support set. Clearly,  $x \in R'_i$ .
- (b)  $x \in R_i, R_j$  such that  $1 \leq i < j \leq m$ . By assumption,  $x$  then does not belong to any other support set and clearly in this case,  $x \in R'_i, R'_j$ .
- (c)  $x \in R_i, R_j$  such that  $i \leq m, j \geq m+1$ . Note that  $x$  then does not belong to any other support set. Now, consider the  $j^{\text{th}}$  column of the matrix  $[I|B]$  and let us call it as  $\mathbf{b}$ . We consider three sub-cases for this situation based on the column weight of  $\mathbf{b}$ .
  - (i) Column weight of  $\mathbf{b}$  is 0. Clearly, then  $x \in R'_i$ .
  - (ii) Column weight of  $\mathbf{b}$  is 1, say  $b_\ell \neq 0$ . Suppose,  $\ell \neq i$ , then  $x \in R'_i, R'_\ell$ . Now if  $\ell = i$ , the element  $x$  need not be present in  $R'_i$ . However, for the purposes of counting  $|\cup_{i=1}^m R'_i|$ , we could replace  $x$  with  $y$  where  $y$  is one of the elements covered only by  $R_j$  (note that the such an element exists by assumption). This works because, if this particular case does occur, we will never again have to seek one of the elements covered only by  $R_j$ . This is because in order for this to happen again it must be true that there exists another element  $x' \in R_i \cap R_j$ , but this is contrary to our assumption that any two support sets have intersection at most 1.
  - (iii) Column weight of  $\mathbf{b}$  is 2 or more, say  $b_{\ell_1}, b_{\ell_2} \neq 0$ . Without loss of generality if assume that  $\ell_1 \neq i$ , then  $x \in S'_{\ell_1}$ .

Thus we see that in all the cases we either do not lose the element  $x$  or there is another unique element  $y$  which can compensate for  $x$  while counting the support cardinality after the linear combinations are taken. Hence we conclude that  $|\cup_{i=1}^m R'_i| \geq |\cup_{i=1}^m R_i|$ .